# Part 4: Regression

Sam Tyner

TBD

# Textbook

These slides are based on the book *OpenIntro Statistics* by David Diez, Christopher Barr, and Mine Çetinkaya-Rundel

The book can be downloaded from
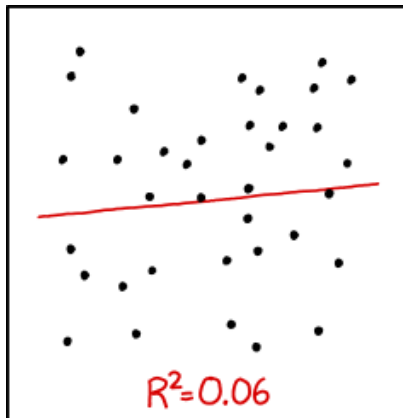https://www.openintro.org/stat/textbook.php

Part 4 Corresponds to Chapters 7, 8 of the text. Sections 4.1-4.3 correspond to chapters 7, 8.1, 8.4 of the text.

# Outline

- Simple Linear Regression (4.1)
- Multiple Regression (4.2)
- Logistic Regression (4.3)

# Section 4.1: Simple Linear Regression

# What is regression?



R²=0.06

REXTHOR, THE DOG-BEARER

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.
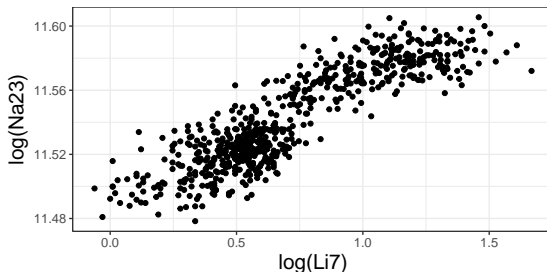
# Formula for a line

$$Y = a \cdot X + b$$

- $Y$: dependent variable (response)
- $X$: independent variable (predictor)
- $a$: slope (for every 1 unit increase in $X$, $Y$ increases by $a$)
- $b$: intercept (when $X = 0$, $Y = b$)
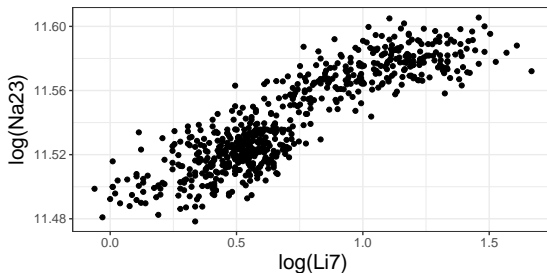- *Deterministic*: knowledge of $a, X, b$ means you know $Y$

# Linear Regression

Why?

- Have two variables, $Y, X$ and we think that the value of $Y$ *depends on* the value of $X$
- Why would we think that? Maybe we have previous knowledge or we looked at a scatterplot of the data
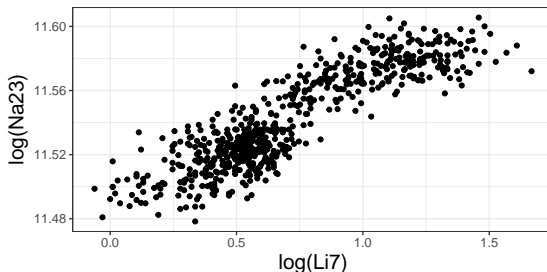
# Linear Regression



Guesstimate the slope and intercept of a line through this data

- $a$: slope = ?
- $b$: intercept = ?

# Linear Regression



Guesstimate the slope and intercept of a line through this data

- *a*: slope ≈ 0.08
- *b*: intercept ≈ 11.48

# Best fit line

The *best fit line* is the equation $Y = a \cdot X + b$ with values $a, b$ that minimize the **sum of squared residuals**. What does that mean?

- Write the best fit line as $\hat{Y} = \beta_0 + \beta_1 \cdot X$.

# Best fit line

The *best fit line* is the equation $Y = a \cdot X + b$ with values $a, b$ that minimize the **sum of squared residuals**. What does that mean?

- Write the best fit line as $\hat{Y} = \beta_0 + \beta_1 \cdot X$.
- $\hat{Y}$ is the predicted value of $Y$ by the best fit line

# Best fit line

The *best fit line* is the equation $Y = a \cdot X + b$ with values $a, b$ that minimize the **sum of squared residuals**. What does that mean?

- Write the best fit line as $\hat{Y} = \beta_0 + \beta_1 \cdot X$.
- $\hat{Y}$ is the predicted value of $Y$ by the best fit line
- **Residual** - what is "left over" after the prediction.

# Best fit line

The *best fit line* is the equation $Y = a \cdot X + b$ with values $a, b$ that minimize the **sum of squared residuals**. What does that mean?

- Write the best fit line as $\hat{Y} = \beta_0 + \beta_1 \cdot X$.
- $\hat{Y}$ is the predicted value of $Y$ by the best fit line
- **Residual** - what is "left over" after the prediction.
- Denote residual for observation $i$ by $e_i = Y_i - \hat{Y}_i$

# Best fit line

The *best fit line* is the equation $Y = a \cdot X + b$ with values $a, b$ that minimize the **sum of squared residuals**. What does that mean?

- Write the best fit line as $\hat{Y} = \beta_0 + \beta_1 \cdot X$.
- $\hat{Y}$ is the predicted value of $Y$ by the best fit line
- **Residual** - what is "left over" after the prediction.
- Denote residual for observation $i$ by $e_i = Y_i - \hat{Y}_i$
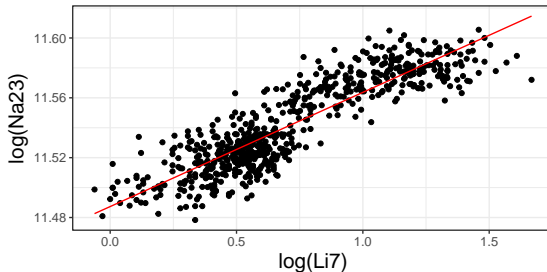- We are minimizing $\sum_{i=1}^{N} e_i^2$

# Calculating the best fit line

- $\beta_1 = \frac{s_y}{s_x} \cdot r$
- $s_y$: standard deviation of the data observations $Y$
- $s_x$: standard deviation of the data observations $X$
- $r$: correlation between the observations $X, Y$ (measure of association between $X, Y$)
- $\beta_0 = \bar{y} - \beta_1 \cdot \bar{x}$

# Do it in R

```r
bfl <- lm(data = glass, log(Na23) ~ log(Li7))
coef(bfl)
```

```
## (Intercept)    log(Li7)
## 11.48732735  0.07632503
```
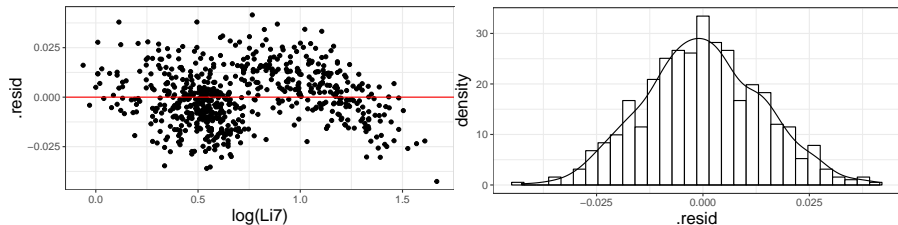
# Residual Plot

Look at the residuals $e$ by the values of $X$:



Want to see a random scatter of points above and below 0

# $R^2$: how well does $X$ explain $Y$?

$R^2$, the **coefficient of determination** defines how much of the variability in $Y$ is explainable by the values of $X$.

$$R^2 = 1 - \frac{Var(e)}{Var(y)}$$

Example:

- $Y = \log(Na23)$. $Var(Y) = 0.00089$
- $e_i = Y_i - \hat{Y}_i = Y_i - (11.487 + 0.0763 \cdot X_i)$. $Var(e) = 0.00019$
- $\frac{Var(e)}{Var(y)} = \frac{0.00019}{0.00089} = 0.2138$
- $R^2 = 1 - \frac{Var(e)}{Var(y)} = 1 - 0.2138 = 0.7862$

78.62% of the variability in $Y$ is explained by the value of $X$.

# Section 4.2: Multiple Regression

# Multiple Regression = Multiple Predictors

Multiple regression is the same general idea as simple linear regression, but instead of one predictor variable, $X$, we have two or more: $X_1, X_2, \ldots, X_p$ where $p$ is the number of predictor variables.

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \cdots + \beta_p \cdot X_p$$

$\beta_0$ (intercept) and $\beta_1, \ldots, \beta_p$ are called *coefficients*

- $\beta_0$ is the value of $\hat{Y}$ when ALL $X$s are 0
- $\beta_k, k \in \{1, 2, \ldots, p\}$ is the amount that $Y$ increases when $X_k$ increases by 1 unit, and *all other X values are held constant*

Why?

# Why Multiple Regression?

- May know that more than 1 variable affects the value of $Y$ (background knowledge)
- More predictors generally means better fit, better predictions

Example: Add log value of Neodymium (common glass additive) to the model

```
blfm2 <- lm(log(Na23) ~ log(Li7) + log(Nd146), data = glass)
blfm2
```

```
##
## Call:
## lm(formula = log(Na23) ~ log(Li7) + log(Nd146), data = glas
##
## Coefficients:
## (Intercept)      log(Li7)     log(Nd146)
##    11.53947       0.05774       -0.03699
```
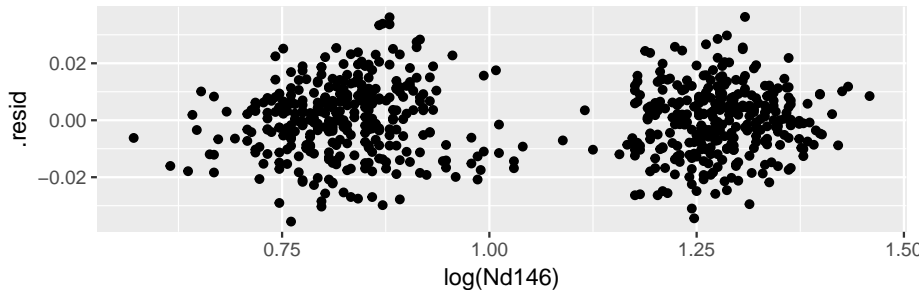
# Multiple Regression Example

Example: $\log(Na23) = \beta_0 + \beta_1 \cdot \log(Li7) + \beta_2 \cdot \log(Nd146)$

Residuals:
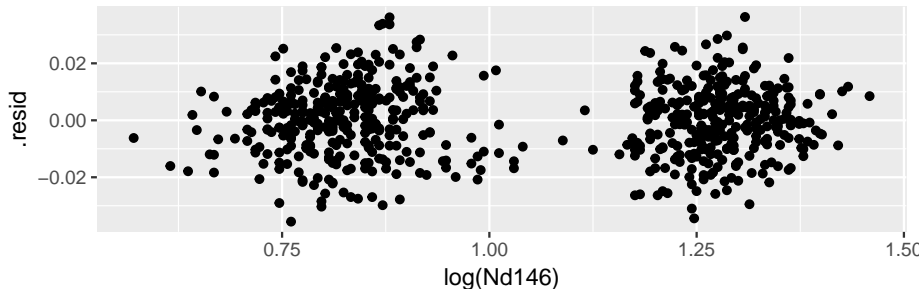


$$R^2 = 1 - \frac{Var(e)}{Var(y)} = 1 - \frac{0.000155}{0.00089} = 0.8258$$

What do you think of this model?

# Multiple Regression Example

Example: $\log(Na23) = \beta_0 + \beta_1 \cdot \log(Li7) + \beta_2 \cdot \log(Nd146)$

Residuals:



$R^2 = 1 - \frac{Var(e)}{Var(y)} = 1 - \frac{0.000155}{0.00089} = 0.8258$

Better fit than the simple linear regression, but we uncovered a new pattern: two distinct groups of residuals

# Another multiple regression

There are 2 manufacturers in the glass data, so we'll add the manufacturer as a variable in the model:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \beta_3 \cdot X_3$$

- $Y = \log(Na23)$, $X_1 = \log(Li7)$, $X_2 = \log(Nd146)$, $X_3 =$ manufacturer.

```
blfm3 <- lm(log(Na23) ~ log(Li7) + log(Nd146) + mfr , data = glass)
blfm3


##
## Call:
## lm(formula = log(Na23) ~ log(Li7) + log(Nd146) + mfr, data = glas
##
## Coefficients:
## (Intercept)     log(Li7)    log(Nd146)         mfrM2
##    11.523955     0.057276     -0.024866      0.006241
```

# Section 4.3: Logistic Regression

# Different types of responses

In linear & multiple regression, we have a *continuous* numerical response variable. However, this is not always the case.

Often, we want to determine whether the response belongs in one of two categories.

Examples:

- Is an email spam or not?
- Is a glass fragment from manufacturer 1 or 2?
- Will a juror say the defendant in a case is guilty or not guilty?

# Logistic regression